

Prediction of 2024 US election ...*

Colin Sihan Yang Lexun Yu Siddharth Gowda

November 3, 2024

We forecast the winner of the 2024 US presidential election using “poll-of-polls” by building a linear model.

1 Introduction

Election result forecasting has become an essential tool for analysts in political science and the public to predict the outcome of democratic process, such as the presidential election in the United States. Traditionally, individual polls have been used as a snapshot of voter sentiment, but they only reflect temporary changes in the performance of contestants, instead of a precise estimation of the election result. As discussed by Pasek (2015) and Blumenthal (2014), the aggregation of multiple polls, or “poll-of-polls,” has become a popular technique to reduce individual survey errors and provide more accurate election forecasts. However, the traditional poll aggregation does not reflect dynamics of an election, especially with differences in quality of polls. A more adaptable model is required to predict the election result based on both polling data and additional variables, such as sample size and pollster credibility.

We build a hybrid election forecasting model following the strategies mentioned by Pasek (2015). As Pasek (2015) described in their article, aggregation involves determining which surveys are worth including, as well as selecting, combining and averaging results from multiple polls to reduce individual biases and errors. Prediction modeling adds other data to the model that predicts election outcomes based on current dynamics. Hybrid models like the Bayesian approach incorporates prior beliefs based on historical data or expert knowledge and new evidence like economic updates to dynamically adjust the forecast as the campaign progresses.

In this paper, we predict the 2024 us election result with the hybrid election forecasting model. We incorporate aggregation by filtering the polls on FiveThirtyEight (2024) by numeric grade that indicates pollster’s reliability, prediction that incorporates social and economic indicators

*Code and data are available at: <https://github.com/yulexun/uselection>.

including unemployment rates and abortion rates, and hybrid approaches that uses Bayesian techniques which combines historical data such as the 2016 election data, allowing for a dynamic prediction of the U.S. presidential election.

The estimand for this research paper is the predicted support percentages for Kamala Harris and Donald Trump. The prediction is based on quantifying various polling factors, including sample size, poll scores, and transparency scores, which are used as predictors.

The results of this model indicate a more stable and accurate forecast compared to traditional aggregation methods alone, [update this ...]

The remainder of this paper is structured as follows: [update this ...]

The data gathering and analysis is done in R (R Core Team 2023) with the following packages: knitr (Xie 2014), tidyverse (Wickham et al. 2019), ggplot2 (Wickham 2016), dplyr (Wickham et al. 2023), arrow (Richardson et al. 2024), here (Müller 2020), gridExtra (Auguie 2017), Matrix (Bates, Maechler, and Jagan 2024), Rstan (Stan Development Team 2024) and lubridate (Grolemund and Wickham 2011).

2 Data

2.1 Overview

For the data we used in this analysis about the polling result for Kamala Harris and Donald Trump in 2024 USA president election.

- **response variable:** `pct`(`pct`: The percentage of the vote or support that the candidate received in the poll)

- **numeric predictor:**

`sample_size`(`sample_size`: The total number of respondents participating in the poll)

`timegap`(the time gap between the poll start date and the real election date i.e `timegap` = `real US election date` - `poll start date`)

`pollscore`(A numeric value representing the score or reliability of the pollster in question)

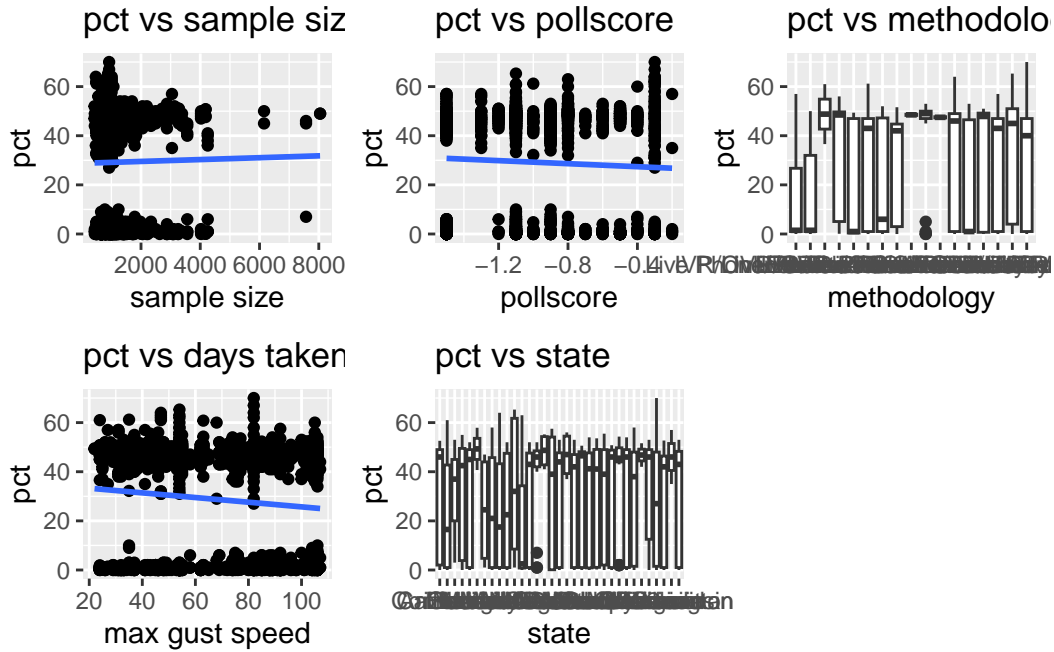
- **categorical predictor** `state`(The U.S. state where the poll was conducted or focused)

`methodology`(The method used to conduct the poll)

2.2 Data Exploration

- `pct` vs `sample_size`: This scatter plot shows the `pct` against the `sample_size`, with a fitted trend line indicating a slight positive relationship. The data points are denser for lower sample sizes, suggesting that smaller sample sizes are more common in the dataset.

Table 1



- pct vs pollscore: This scatter plot illustrates the pct against pollscore. The fitted trend line suggests a weak negative relationship between pollscore and pct. The points are scattered without a strong linear pattern.
- pct vs methodology: A boxplot comparing pct for different polling methodologies. The pct distribution varies across methodologies, with some showing greater spread or median differences. This suggests that the polling methodology may influence pct outcomes.
- pct vs days taken from election: A scatter plot displaying pct versus the number of days before the election. The trend line indicates a slight negative relationship, suggesting that as the election date approaches, pct may decrease slightly.
- pct vs state: A boxplot depicting pct across different states. The pct distribution varies by state, with some states showing wider variability or different median values, implying state-specific effects on pct.

Figure 1 The pairs plot displays scatter plots of four numeric variables (`pct`, `sample_size`, `pollscore`, and `days_taken_from_election`) to visualize their relationships. The data shows clustering, particularly in `pct` versus `sample_size`, suggesting potential heteroscedasticity. The `sample_size` variable is skewed towards lower values, while `pollscore` and `days_taken_from_election` have a more even spread, though `pollscore` shows central clustering. No strong linear relationships are immediately apparent between the variables, indicating that correlations are likely weak.

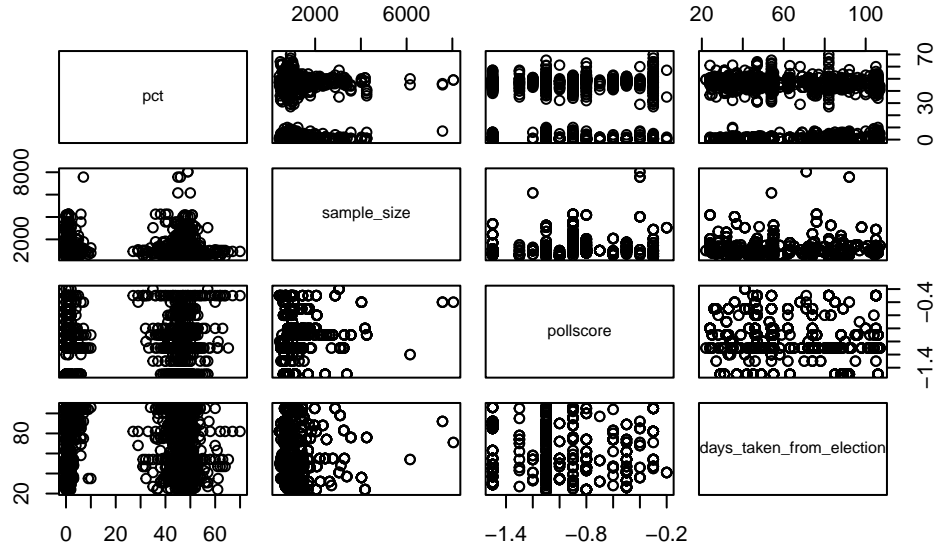


Figure 1

2.3 Measurement

In this dataset, each row represents a polling question that records the variables of interest. Each entry allows us to explore the real-world relationships between polling factors and the support percentage (`pct`) for the candidates Kamala Harris and Donald Trump. This dataset enables an analysis of how various polling characteristics influence the reported support levels for the candidates we are focused.

2.4 Clean Data

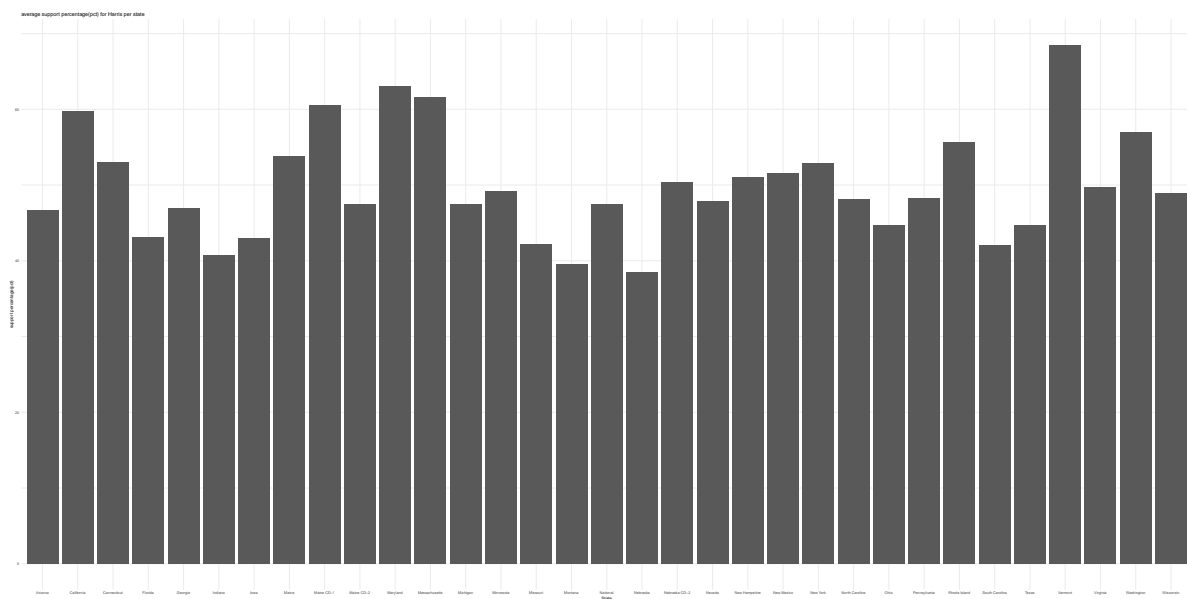
Table 2 The data cleaning process involves several steps to ensure the quality and relevance of the polling data. First, we filter the dataset to retain only poll results with a numeric grade of 2.7 or higher, indicating that the polls are considered reliable. Next, we address missing values in the state attribute: polls with NA in the state column are considered national polls.

We then create a new attribute, `days_taken_from_election`, which represents the time gap between the poll's start date and the actual U.S. election date. Additionally, we filter the dataset to include only polls conducted after July 21, 2024, the date when Kamala Harris declared her candidacy. Finally, we remove any remaining rows that contain missing values to ensure a clean dataset.

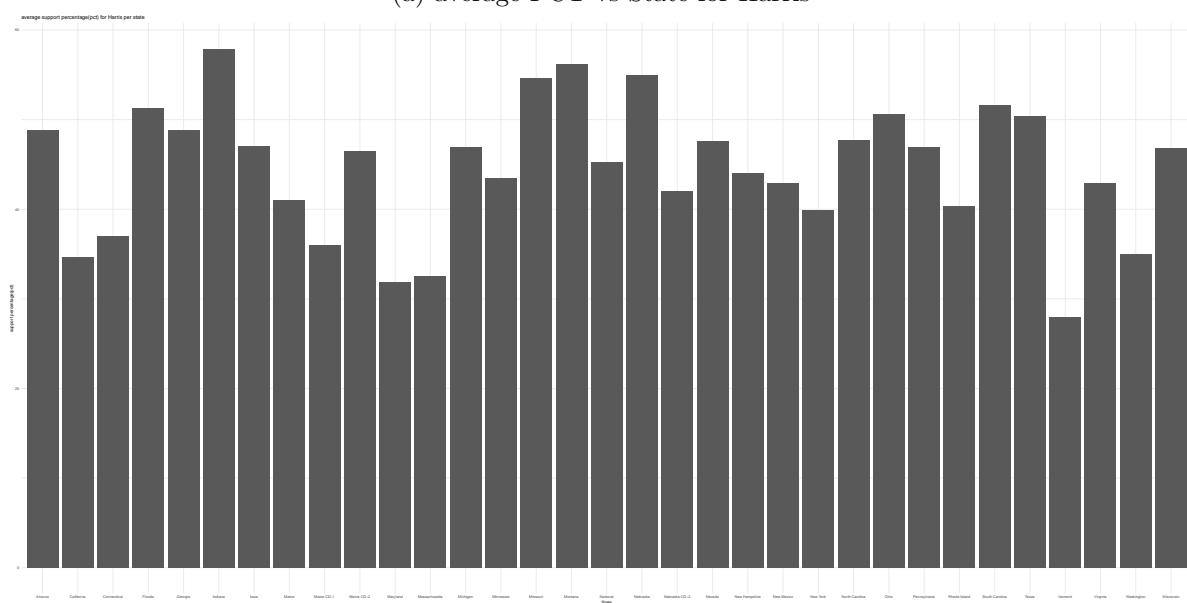
Table 2: Sample of cleaned US election data

pct	sample_size	pollscore	days_taken_from_election	state	methodology	candidate_name
47.6	4180	-0.8	24	National	Online Ad	Kamala Harris
50.7	4180	-0.8	24	National	Online Ad	Donald Trump
0.8	4180	-0.8	24	National	Online Ad	Jill Stein
0.1	4180	-0.8	24	National	Online Ad	Chase Oliver
0.1	4180	-0.8	24	National	Online Ad	Cornel West
48.1	4180	-0.8	24	National	Online Ad	Kamala Harris

2.5 Basic Statistics Summary for Data



(a) average PCT vs State for Harris



(b) PCT vs State for Trump

Figure 2: the average PCT vs State for Harris and Trump

Figure 2a The histogram displays the average support percentage for Kamala Harris across different U.S. states. The data indicates that support for Harris varies significantly across states. Notable observations include relatively high average support in states such as California and New York, which are known for being more Democratic-leaning. On the other hand, there are states with lower average support percentages, particularly in more traditionally Republican or swing states. The distribution suggests regional variations in support, with some states showing consistent backing for Harris while others indicate a weaker performance.

Figure 2b The second histogram shows the average support percentage for Donald Trump across various states. It highlights substantial support in states such as Florida and Texas, which align with historical trends of strong Republican support. Trump’s average support appears robust in many midwestern and southern states, which are known for their conservative voter base. However, in more liberal-leaning states such as California and New York, the average support is lower, reflecting these states’ tendency to lean Democratic.

3 Model

The goal of our modeling strategy is twofold. Firstly, we aim to accurately predict the support percentage (PCT) for Harris and Trump based on relevant poll data and key influencing factors. Secondly, we seek to evaluate the efficacy of different modeling approaches—from simple linear regression (SLR) to multiple linear regression (MLR) and Bayesian hierarchical models—to understand their predictive capabilities and assess the underlying relationships between variables. By comparing these models, we can determine which approach provides the most robust and reliable predictions, while considering the variability and potential uncertainty in the data.

3.1 Model set-up

The Bayesian model is implemented in R (R Core Team 2023) using the `rstanarm` package as described by Brilleman et al. (2018). The model is run with the following specifications:

- Formula: $\text{pct} \sim \text{pollscore} + \text{days taken from election} + \text{sample size} + (1|\text{methodology}) + (1|\text{state})$
- Priors: $\text{Normal}(0, 2.5)$ for all coefficients and intercept, $\text{Exponential}(1)$ for σ
- Settings: Seed = 123, Cores = 4, Adapt delta = 0.95

We run the model in R (R Core Team 2023) using the `rstanarm` package of Brilleman et al. (2018). We use the default priors from `rstanarm`. us

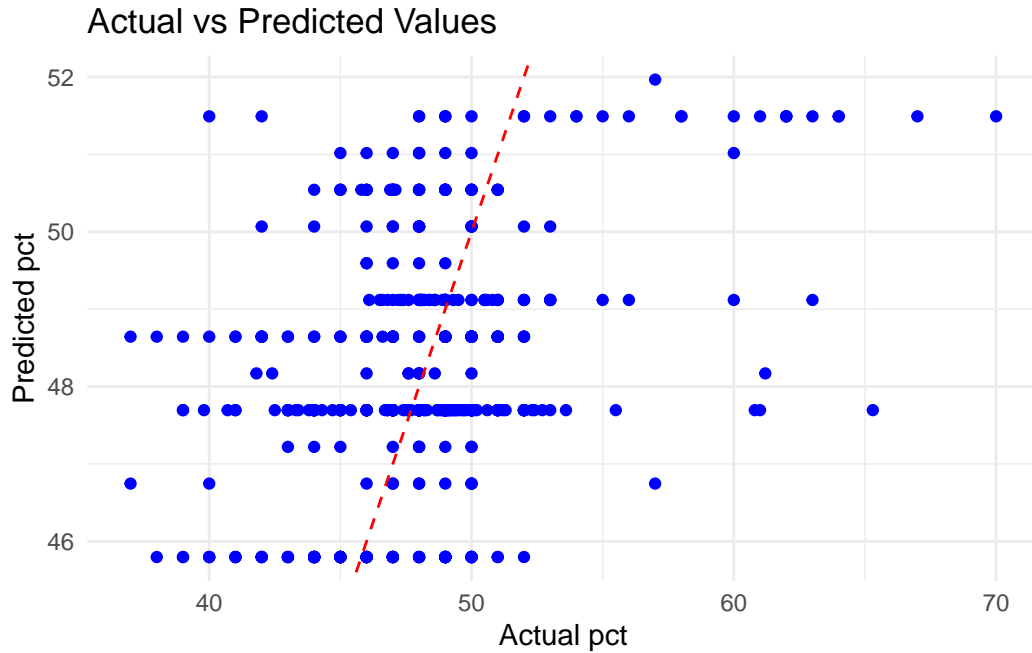


Figure 3

3.2 Basic Model

Figure 3 in this SLR model, the response variable is pct and the only one predictor is the pollscore. The primary concern lies in the evident dispersion of data points, which are widely spread and do not cluster closely around the line of perfect prediction (the dashed red line). This suggests that while pollscore may have some predictive capability, it does not adequately explain the variability in pct. The observed inconsistencies between actual and predicted values indicate that the relationship between pct and pollscore is likely not sufficiently captured by a linear model with just one predictor.

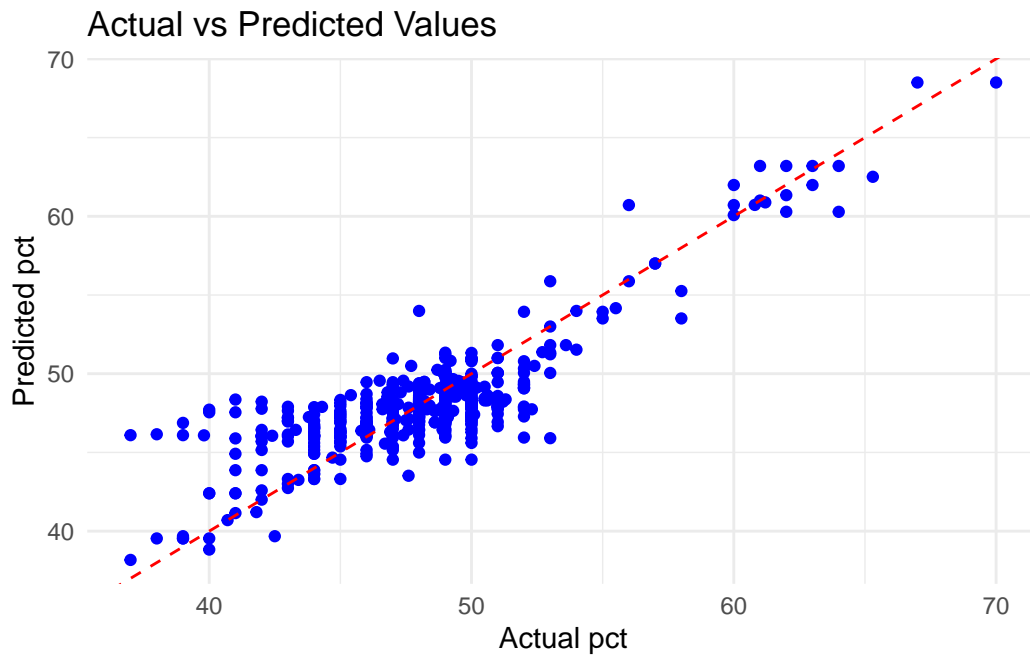
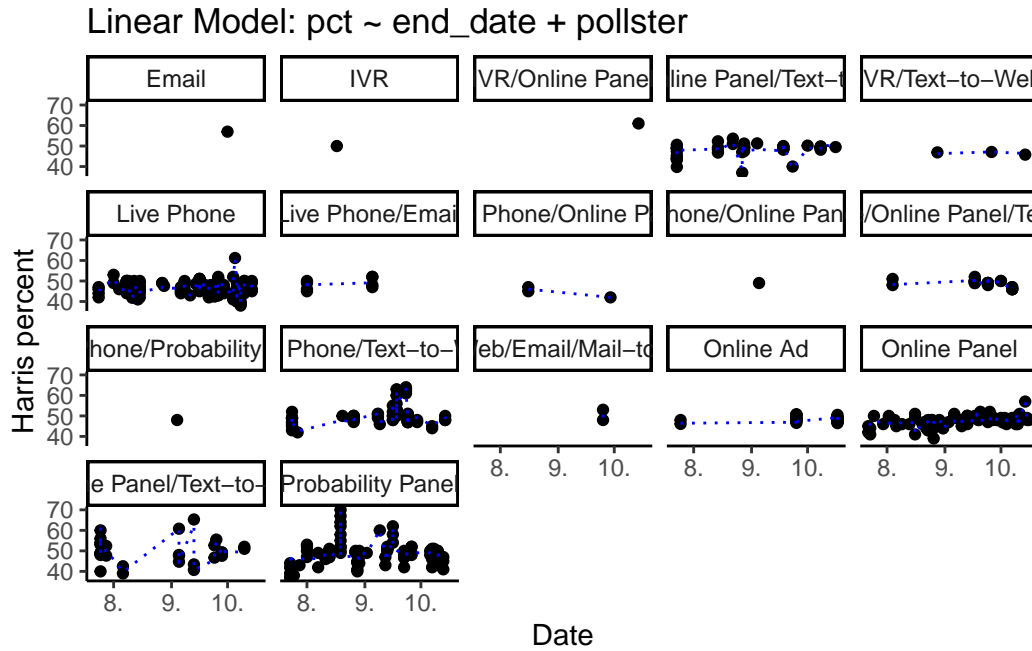


Figure 4

Figure 4, in this MLR model, we predicts pct using pollscore, days_taken_from_election, methodology, sample_size, and state as predictors. The overall distribution of points suggests

that the MLR model has captured a substantial portion of the variance in pct. The majority of the data points align relatively well with the red dashed line, particularly within the middle range of pct values (approximately between 40 and 60). This alignment indicates that the model performs reasonably well in this range, with predicted values correlating strongly with the actual observed outcomes.

3.3 Bayesian Model

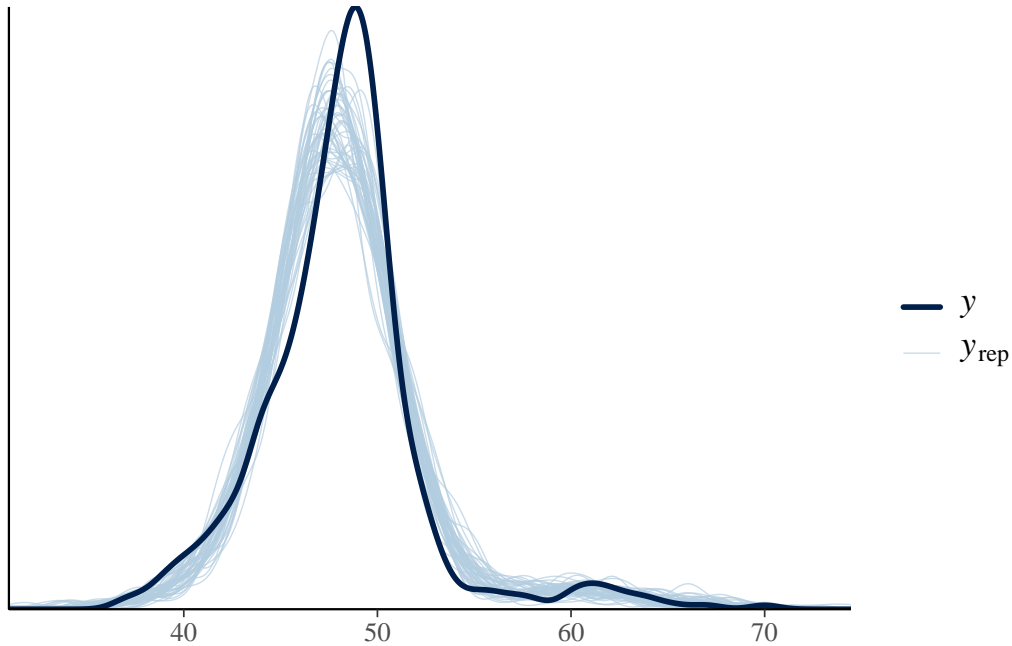


Figure 5

3.3.1 Posterior Predictive Checks

Figure 5 This plot displays the results of a posterior predictive check (PPC) for a Bayesian model using the `pp_check()` function.

The observed data (y) follows a distinct, symmetrical distribution centered around a specific value, suggesting a well-defined peak with tapering on both sides. The replicated distributions (y_{rep}), shown as multiple thin lines, generally align with the shape of the observed distribution, indicating that the Bayesian model has captured the main characteristics of the data. However, the variability in the y_{rep} lines highlights the degree of uncertainty inherent in the model's predictive capability.

The fact that the y_{rep} curves closely match the overall pattern of the actual y suggests that the model performs reasonably well in replicating the observed data. Minor discrepancies or deviations between the y and y_{rep} might imply areas where the model could be fine-tuned or adjusted to improve accuracy. Overall, this PPC indicates that the model provides a decent fit to the data, with some variability accounted for in the predictive samples.

3.3.2 Train Test Validation

```
# Split data into training and test sets
set.seed(123)
train_indices <- sample(seq_len(nrow(just_harris_data)), size = 0.7 * nrow(just_harris_data))
train_data <- just_harris_data[train_indices, ]
test_data <- just_harris_data[-train_indices, ]

# Fit the model on the training data
bayesian_model_train <- stan_glmer(
  formula = formula,
  data = train_data,
  family = gaussian(),
  prior = priors,
  prior_intercept = priors,
  seed = 123,
  cores = 4,
  adapt_delta = 0.95
)

# Predict on the test data and check performance
predictions <- posterior_predict(bayesian_model_train, newdata = test_data)

ss_total <- sum((test_data$pct - mean(test_data$pct))^2)
ss_residual <- sum((test_data$pct - colMeans(predictions))^2)
r_squared <- 1 - (ss_residual / ss_total)
print(paste("R-squared:", r_squared))
```

```
[1] "R-squared: 0.700566347978073"
```

```
# "R-squared: 0.700563202695799"
```

$$SS_{\text{total}} = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (1)$$

$$SS_{\text{residual}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (2)$$

$$R^2 = 1 - \frac{SS_{\text{residual}}}{SS_{\text{total}}} \quad (3)$$

$$\text{R-squared} = 0.700563202695799 \quad (4)$$

Where:

- y_i represents the actual values,
- \bar{y} is the mean of the actual values,
- \hat{y}_i are the predicted values,
- SS_{total} is the total sum of squares,
- SS_{residual} is the sum of squared residuals.

we could see that the R squared value equal to the 0.700563202695799 suggesting that the model captures a substantial portion of the variability in the data and demonstrates that the model performs well in explaining the variability of the pct in the test data, reflecting strong predictive capabilities.

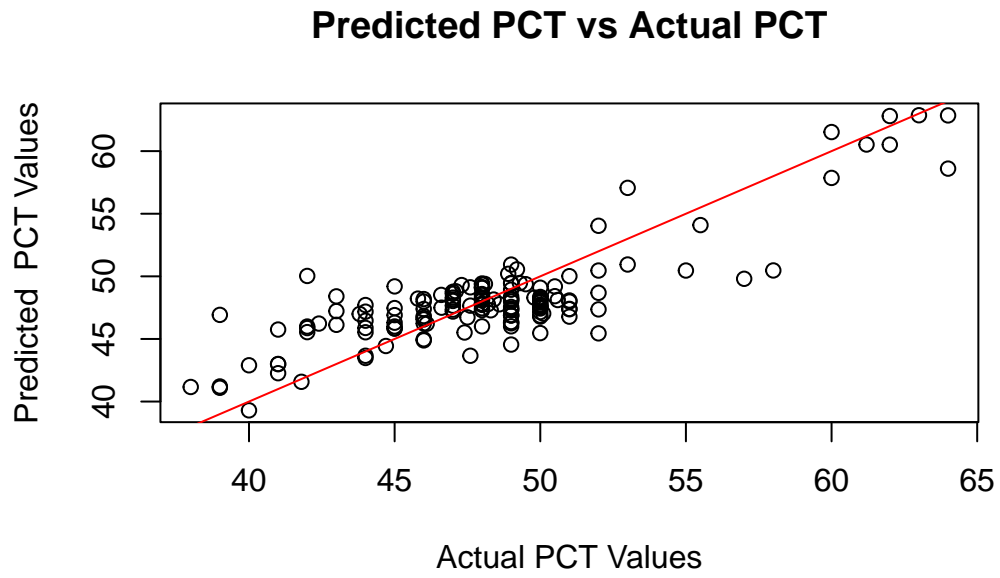


Figure 6

by visual checking actual pct and predicted pct plot Figure 6 the points are plotted against a 45-degree line, which represents the ideal scenario where predicted values match the actual values perfectly.

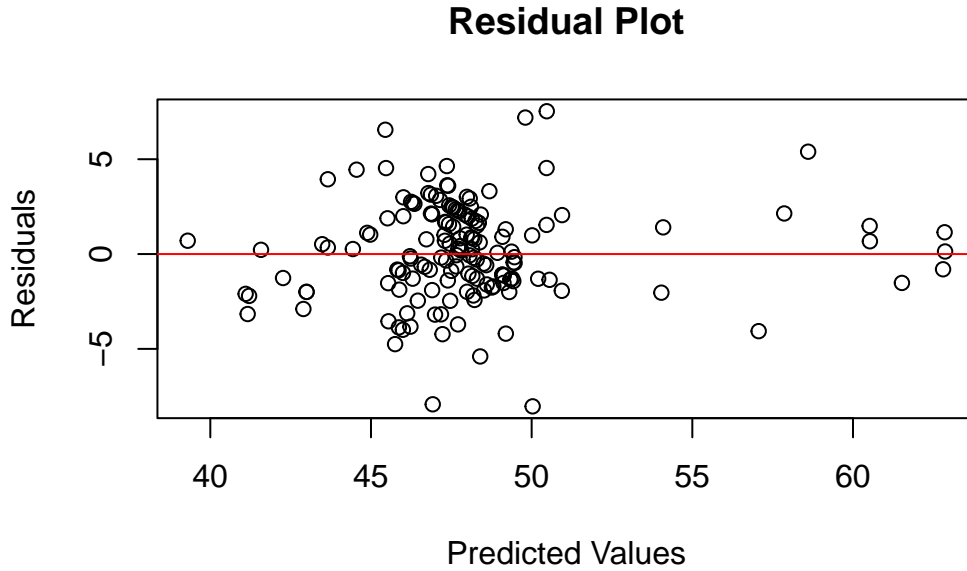


Figure 7

and the residual plot Figure 7 are fairly centered around zero with no major trend, suggesting that the model is not heavily biased in its predictions.

In conclusion, the visual checks from the predicted vs. actual plot and the residual plot, there is no strong evidence that the Bayesian model is overfitting. It appears to generalize well to the data it was trained on without showing signs of capturing noise or irrelevant patterns.

3.3.3 Model justification

To predict the support percentage (PCT) for Harris and Trump, we employed a comprehensive modeling strategy involving Simple Linear Regression (SLR), Multiple Linear Regression (MLR), and a Bayesian hierarchical model. The SLR model provided a foundational analysis of the relationship between PCT and pollscore, highlighting its limited predictive power due to a lack of complexity. The MLR model improved on this by incorporating additional predictors such as days taken from the election, sample size, methodology, and state, enhancing its predictive accuracy and capturing interactions among variables. The Bayesian hierarchical model further refined our approach, incorporating prior knowledge and accounting for variability at group levels (e.g., methodology and state) through random intercepts. This model provided robust uncertainty quantification through credible intervals and demonstrated strong predictive performance with an R-squared value around 0.70, indicating it effectively explained data variability. By leveraging the strengths of these models, particularly the Bayesian approach's

interpretability and robustness, we achieved reliable predictions and a deeper understanding of the factors influencing support percentages.

4 Result

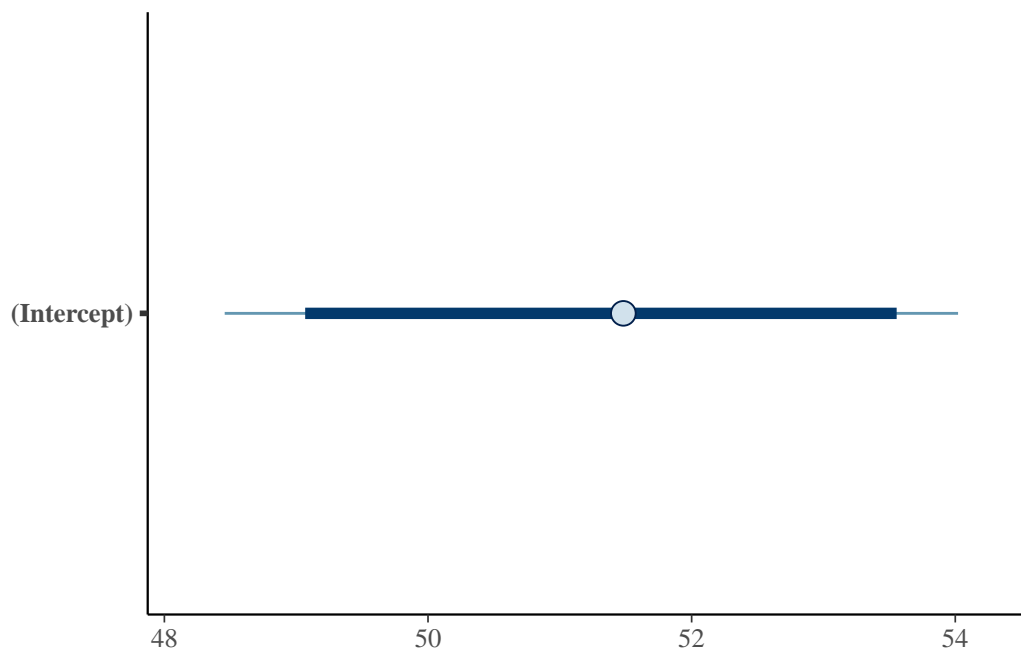


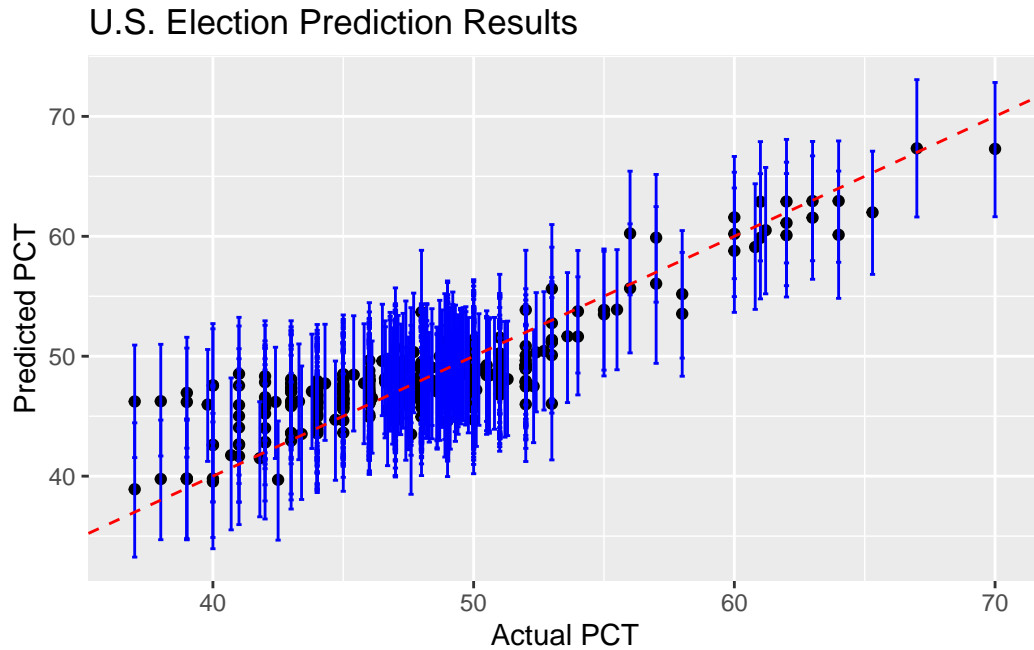
Figure 8

Figure 8 This plot displays the estimated posterior distribution for the intercept parameter in the Bayesian model. The intercept represents the baseline level of support, expressed as a percentage, for Kamala Harris in the U.S. election data used in the analysis. The plot features a point estimate (depicted by the central dot) that indicates the mean or median of the posterior distribution for the intercept, and a horizontal line showing the 95% credible interval, which signifies the range within which the true value of the intercept is likely to fall with 95% probability.

The range of this credible interval spans from approximately 50 to 54 percent, suggesting that the model estimates the baseline level of support for Kamala Harris to be around this range. The relatively narrow width of the interval implies a certain degree of confidence in the estimate, indicating that the data used in the model provided a clear signal for the intercept's value.

In conclusion, the estimated intercept, which represents the baseline percentage of support for Kamala Harris in the U.S. election data analyzed, is centered around 52%, with a credible

interval spanning approximately from 50% to 54%. This suggests that, according to the Bayesian model, the underlying support for Kamala Harris.



	Actual_PCT	Predicted_PCT	Lower_CI	Upper_CI
1	47.6	49.12850	44.44701	54.04904
2	48.1	49.13812	44.34297	53.83691
3	48.6	47.74757	42.95339	52.43367
4	49.3	47.71647	42.90265	52.32287
5	48.1	48.02783	43.50669	52.78940
6	48.4	47.98070	43.11644	52.70039

5 Discussion

5.1 First discussion point

If my paper were 10 pages, then should be be at least 2.5 pages. The discussion is a chance to show off what you know and what you learnt from all this.

5.2 Second discussion point

Please don't use these as sub-heading labels - change them to be what your point actually is.

5.3 Third discussion point

5.4 Weaknesses and next steps

Weaknesses and next steps should also be included.

Appendix

A Additional data details

B Model details

B.1 Posterior predictive check

In [?@fig-ppcheckandposteriorvsprior-1](#) we implement a posterior predictive check. This shows...

In [?@fig-ppcheckandposteriorvsprior-2](#) we compare the posterior with the prior. This shows...

B.2 Diagnostics

[?@fig-stanareyouokay-1](#) is a trace plot. It shows... This suggests...

[?@fig-stanareyouokay-2](#) is a Rhat plot. It shows... This suggests...

C Methodology of YouGov

YouGov’s methodology documentations are separated in two articles. The article by Bailey and Rivers (2024) documents the methodology of the 2024 election projection, while the webpage on YouGov (n.d.) documents the general methodology of YouGov’s prediction.

C.1 Population, Frame, and Sample

As Bailey and Rivers (2024) stated, the population covered by YouGov’s MRP model is everyone in the national voter file, whether or not they belong to YouGov’s panel. The national voter files are digital database built by commercial organizations with public government records of voters, as explained by DeSilver (2018). Voter files indicates whether someone voted in a given election, thus YouGov’s population covers all voters in previous US elections.

YouGov’s sampling frame consists of its online panel members. These members are part of the SAY24 project, a collaboration between Stanford, Arizona State, and Yale Universities, as stated by Bailey and Rivers (2024). YouGov collect information on respondents when they join their panel before they are invited to participate in the survey.

YouGov select the sample from the sampling frame based on their ability to match characteristics of the population of interest. YouGov interviews nearly 100,000 people in the first set of estimates. For the second set of estimates, YouGov didn't just start over with a new sample. They took the initial data from August and September and updated it with responses from more than 20,000 additional registered voters who were re-interviewed in late September and early October.

C.2 Sample Recruitment

Panelists are recruited through various online channels, including advertisements and partnerships with websites (YouGov n.d.). They must provide demographic details upon joining, which helps in selecting representative samples for each survey. When respondents complete a survey, they are awarded points that can be exchanged for money.

C.3 Sampling Approach and Trade-offs

YouGov uses non-probability sampling due to the compensation, an approach where not every individual has an equal chance of selection (YouGov n.d.). This method allows quick and cost-effective data collection. However, as YouGov (n.d.) writes the panelists must have an internet connection to participate. YouGov state that there is 95% of us population with internet access, thus the sample may be less representative of certain hard-to-reach populations, such as individuals with very slow internet access or without internet access.

C.4 Non-response Handling

YouGov apply statistical weighting to adjust for the differences between the sample and target population. The weight is based on demographic characteristics such as age, gender, race and presidential vote (YouGov n.d.). Additionally, quality control measures exclude unreliable responses to improve data accuracy. The respondents are offered a small incentive to decrease the non-response and increase participation.

C.5 Strengths and Weaknesses of the Questionnaire

YouGov's surveys are conducted online, which is very efficient for the respondents, and responses are weighted to enhance representativeness. The pollster can recruit a large amount of panelists because of the online format. Combining with online tracking technologies, the metadata provided by their panelists can be verified easily.

As a non-probability sample, it might miss certain demographic groups not covered by the online population. While weighting improves accuracy, it cannot fully substitute the randomization found in probability sampling. Additionally, the categories in the survey is oversimplified with bias. For instance, in the poll result published by YouGov, gender is divided into Male and Female. Race is divided into White, Black, Hispanic and Other. This indicates a lack of representation.

D FiftyEight Licenses

FiftyEight's data sets are used and modified by us under the [Creative Commons Attribution 4.0 International License](#).

References

- Auguie, Baptiste. 2017. *gridExtra: Miscellaneous Functions for "Grid" Graphics*. <https://CRAN.R-project.org/package=gridExtra>.
- Bailey, Delia, and Douglas Rivers. 2024. "How YouGov's MRP Model Works for the 2024 U.S. Presidential and Congressional Elections." *YouGov*. <https://today.yougov.com/politics/articles/50587-how-yougov-mrp-model-works-2024-presidential-congressional-elections-polling-methodology>.
- Bates, Douglas, Martin Maechler, and Mikael Jagan. 2024. *Matrix: Sparse and Dense Matrix Classes and Methods*. <https://CRAN.R-project.org/package=Matrix>.
- Blumenthal, Mark. 2014. "Polls, Forecasts, and Aggregators." *PS: Political Science and Politics* 47 (2): 297–300. <http://www.jstor.org/stable/43284537>.
- Brilleman, SL, MJ Crowther, M Moreno-Betancur, J Bueros Novik, and R Wolfe. 2018. "Joint Longitudinal and Time-to-Event Models via Stan." https://github.com/stan-dev/stancon_talks/.
- DeSilver, Drew. 2018. "Q&A: The Growing Use of 'Voter Files' in Studying the U.S. Electorate." *Pew Research Center*. <https://www.pewresearch.org/short-reads/2018/02/15/voter-files-study-qa/>.
- FiveThirtyEight. 2024. "Our Data." *FiveThirtyEight*. <https://data.fivethirtyeight.com>.
- Grolemund, Garrett, and Hadley Wickham. 2011. "Dates and Times Made Easy with lubridate." *Journal of Statistical Software* 40 (3): 1–25. <https://www.jstatsoft.org/v40/i03/>.
- Müller, Kirill. 2020. *Here: A Simpler Way to Find Your Files*. <https://CRAN.R-project.org/package=here>.
- Pasek, Josh. 2015. "THE POLLS—REVIEW: PREDICTING ELECTIONS: CONSIDERING TOOLS TO POOL THE POLLS." *The Public Opinion Quarterly* 79 (2): 594–619. <http://www.jstor.org/stable/24546379>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Richardson, Neal, Ian Cook, Nic Crane, Dewey Dunnington, Romain François, Jonathan Keane, Dragoş Moldovan-Grünfeld, Jeroen Ooms, Jacob Wujciak-Jens, and Apache Arrow. 2024. *Arrow: Integration to 'Apache' 'Arrow'*. <https://CRAN.R-project.org/package=arrow>.
- Stan Development Team. 2024. "RStan: The R Interface to Stan." <https://mc-stan.org/>.
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. "Welcome to the tidyverse." *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- Wickham, Hadley, Romain François, Lionel Henry, Kirill Müller, and Davis Vaughan. 2023. *Dplyr: A Grammar of Data Manipulation*. <https://CRAN.R-project.org/package=dplyr>.
- Xie, Yihui. 2014. "Knitr: A Comprehensive Tool for Reproducible Research in R." In *Implementing Reproducible Computational Research*, edited by Victoria Stodden, Friedrich

Leisch, and Roger D. Peng. Chapman; Hall/CRC.
YouGov. n.d. “Methodology.” Accessed October 31, 2024. [https://today.yougov.com/about/
panel-methodology](https://today.yougov.com/about/panel-methodology).